# DIABETES PREDICATION ASSESSMENT

TANYA MITTAL (DATA ANALYST INTERN)

# INTRODUCTION

The provided dataset by psyliq comprises detailed information about diabetes patients, which includes patient ID, age, gender, body mass index, blood pressure and few more columns related to it. This dataset comprises approximately 100,000 rows and 11 columns. To extract relevant insights from the dataset, the analysis has been performed using the MYSQL tool, which focuses on extracting meaningful insights. This tool provides a comprehensive and flexible environment for working with large datasets, enabling analysts to query and manipulate the data in various ways.
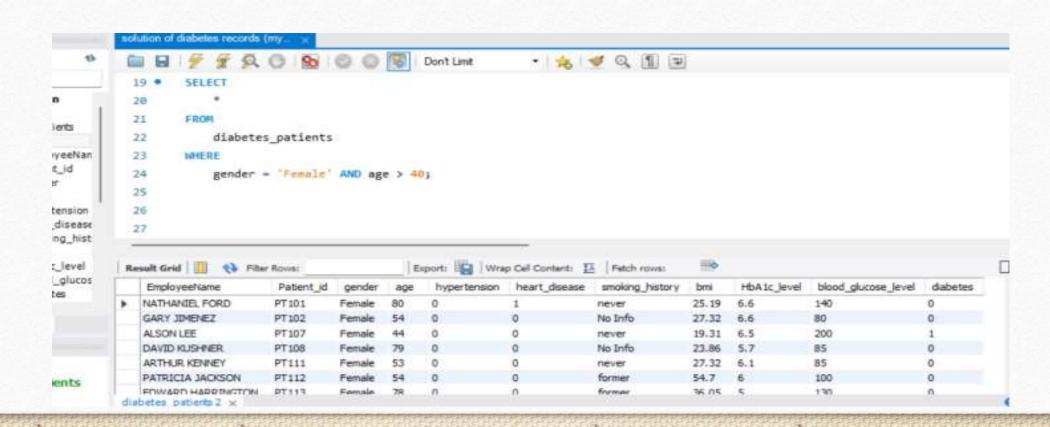
# Retrieve the Patient_id and ages of all patients.

# Select all female patients who are older than 40.

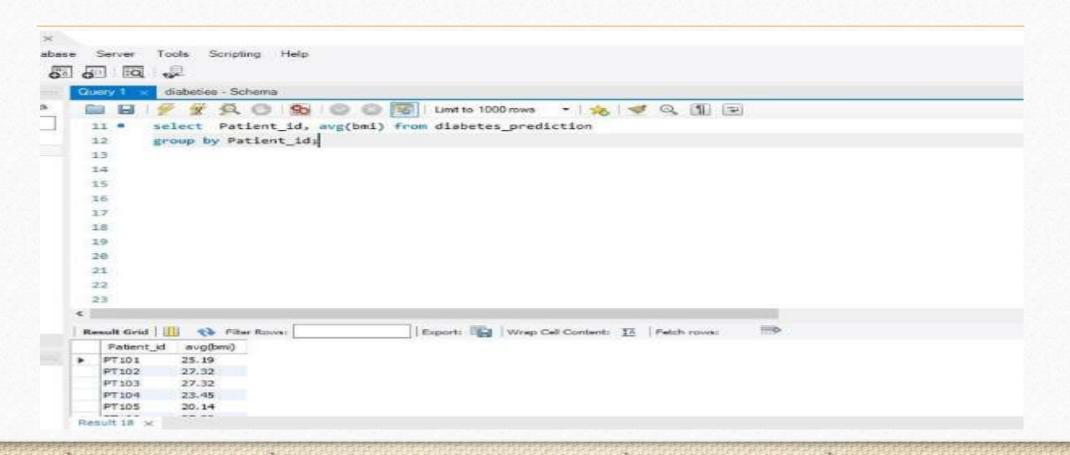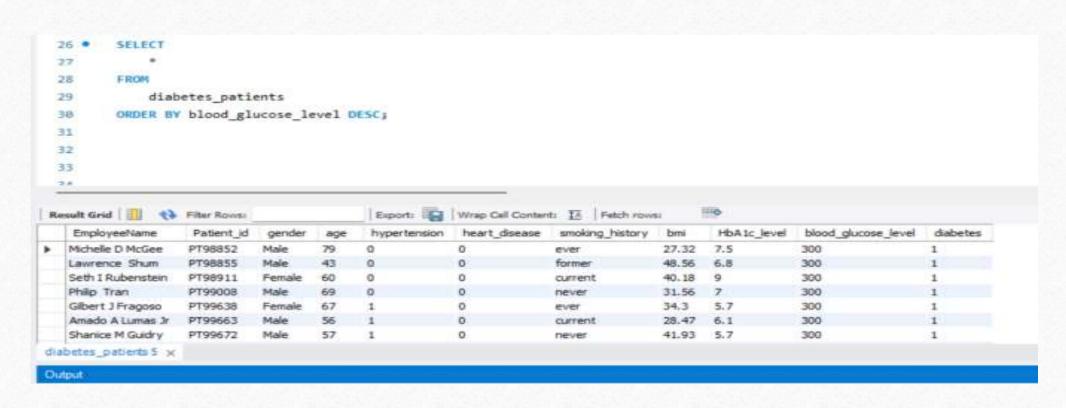# Calculate the average BMI of patients.

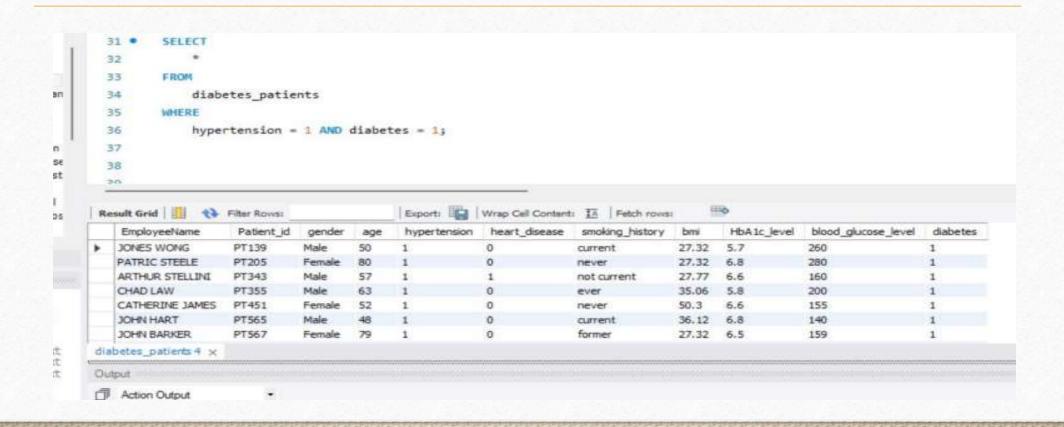# List patients in descending order of blood glucose levels.

```
26 •     SELECT
27           *
28       FROM
29           diabetes_patients
30       ORDER BY blood_glucose_level DESC;
31
32
33
```

**Result Grid** | Filter Rows: | Export: | Wrap Cell Content: | Fetch rows:

| | EmployeeName | Patient_id | gender | age | hypertension | heart_disease | smoking_history | bmi | HbA1c_level | blood_glucose_level | diabetes |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ► | Michelle D McGee | PT98852 | Male | 79 | 0 | 0 | ever | 27.32 | 7.5 | 300 | 1 |
| | Lawrence Shum | PT98855 | Male | 43 | 0 | 0 | former | 48.56 | 6.8 | 300 | 1 |
| | Seth I Rubenstein | PT98911 | Female | 60 | 0 | 0 | current | 40.18 | 9 | 300 | 1 |
| | Philip Tran | PT99008 | Male | 69 | 0 | 0 | never | 31.56 | 7 | 300 | 1 |
| | Gilbert J Fragoso | PT99638 | Female | 67 | 1 | 0 | ever | 34.3 | 5.7 | 300 | 1 |
| | Amado A Lumas Jr | PT99663 | Male | 56 | 1 | 0 | current | 28.47 | 6.1 | 300 | 1 |
| | Shanice M Guidry | PT99672 | Male | 57 | 1 | 0 | never | 41.93 | 5.7 | 300 | 1 |

diabetes_patients 5 ✕

Output

# Find patients who have hypertension and diabetes.

# Determine the number of patients with heart disease.

```
43 •   SELECT
44         COUNT(*)
45     FROM
46         diabetes_patients
47     WHERE
48         heart_disease = 1;
49
```
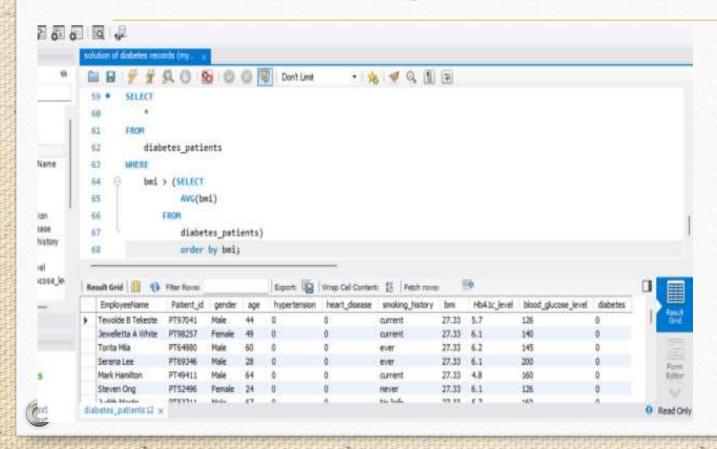
| count(*) |
| --- |
| 3942 |

# Group patients by smoking history and count how many smokers and non-smokers there are.

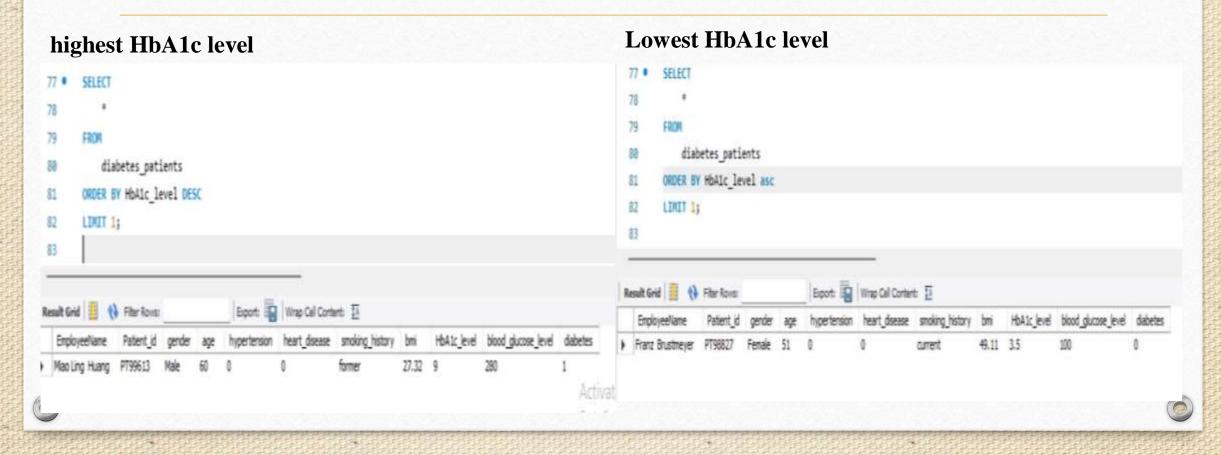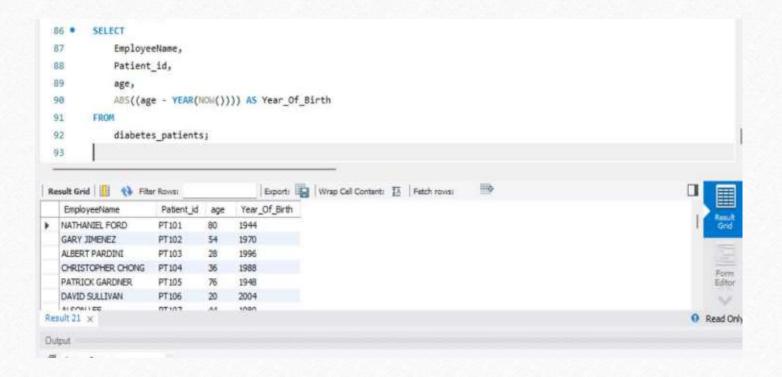# Retrieve the Patient_ids of patients who have a BMI greater than the average BMI.

# Find the patient with the highest HbA1c level and the patient with the lowest HbA1clevel.

**highest HbA1c level**

```
77  •  SELECT
78         *
79      FROM
80         diabetes_patients
81      ORDER BY HbA1c_level DESC
82      LIMIT 1;
83
```

| EmployeeName | Patient_id | gender | age | hypertension | heart_disease | smoking_history | bmi | HbA1c_level | blood_glucose_level | diabetes |
|---|---|---|---|---|---|---|---|---|---|---|
| Mao Ling Huang | PT99613 | Male | 60 | 0 | 0 | former | 27.32 | 9 | 280 | 1 |

**Lowest HbA1c level**

```
77  •  SELECT
78         *
79      FROM
80         diabetes_patients
81      ORDER BY HbA1c_level asc
82      LIMIT 1;
83
```

| EmployeeName | Patient_id | gender | age | hypertension | heart_disease | smoking_history | bmi | HbA1c_level | blood_glucose_level | diabetes |
|---|---|---|---|---|---|---|---|---|---|---|
| Franz Brustmeyer | PT98827 | Female | 51 | 0 | 0 | current | 49.11 | 3.5 | 100 | 0 |

# Calculate the age of patients in years (assuming the current date as of now).

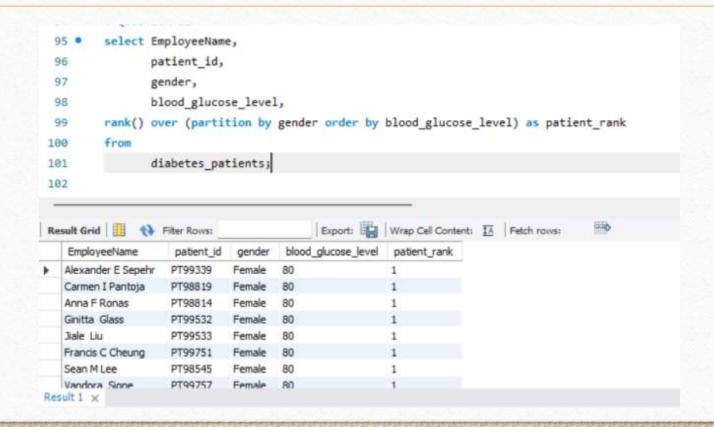# Rank patients by blood glucose level within each gender group.

# Update the smoking history of patients who are older than 50 to "Ex-smoker."
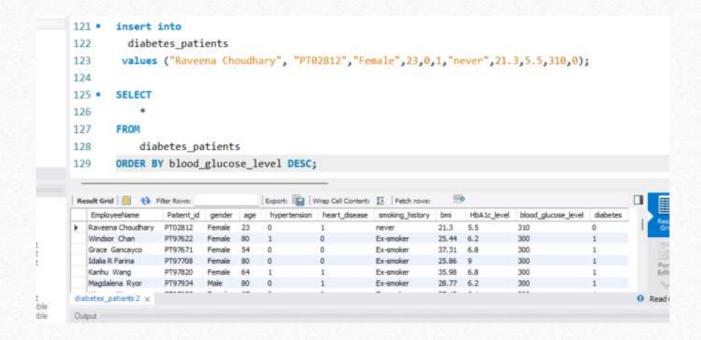
# Insert a new patient into the database with sample data.

# Delete all patients with heart disease from the database.

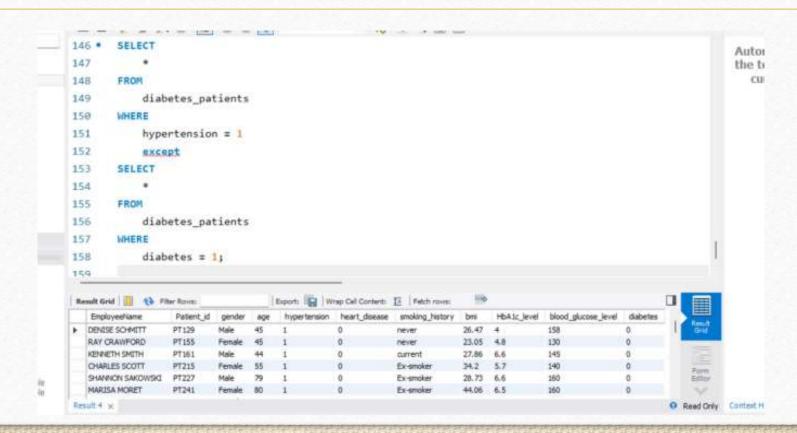# Find patients who have hypertension but not diabetes using the EXCEPT operator.

# Define a unique constraint on the "patient_id" column to ensure its values are unique.

**Before**

| Column Name | Datatype | PK | NN | UQ |
|---|---|:---:|:---:|:---:|
| ◇ EmployeeName | VARCHAR(100) | ☐ | ☑ | ☐ |
| 🔑 Patient_id | VARCHAR(45) | ☑ | ☑ | ☐ |
| ◇ Gender | VARCHAR(20) | ☐ | ☐ | ☐ |

**Query**

```
ALTER TABLE patient_data ADD CONSTRAINT UNIQUE (patient_id);
```

| ✅ | 16 | 19:56:55 | ALTER TABLE patient_data ADD CONSTRAINT UNIQUE (patient_id) | 0 row(s) affected Records: 0 Duplicates: 0 Warnings: 0 |
|---|---|---|---|---|

**After**

| Column Name | Datatype | PK | NN | UQ |
|---|---|:---:|:---:|:---:|
| ◇ EmployeeName | VARCHAR(100) | ☐ | ☑ | ☐ |
| 🔑 Patient_id | VARCHAR(45) | ☑ | ☑ | ☑ |
| ◇ Gender | VARCHAR(20) | ☐ | ☐ | ☐ |

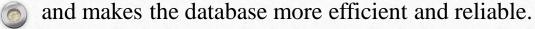# Create a view that displays the Patient_ids, ages, and BMI of patients.

# Suggest improvements in the database schema to reduce data redundancy and improve data integrity.

- Normalize your database schema to minimize redundancy. This involves breaking down large tables into smaller ones and establishing relationships between them through foreign keys. This ensures that each piece of data is stored only once, reducing redundancy.
- Implement foreign keys to enforce referential integrity between related tables. This ensures that relationships between tables are maintained, preventing orphaned records and ensuring data consistency.
- Define constraints and triggers to enforce business rules and maintain data integrity at the database level. Constraints such as NOT NULL, CHECK, and DEFAULT can help ensure that only valid data is entered into the database.
- Utilize unique constraints on columns where appropriate. This ensures that each value in a particular column is unique, preventing duplicates and improving data integrity.
- Determine the primary key for each entity, which uniquely identifies each record in the table.
- Break down a larger table into smaller tables and establish relationships between them. This reduces redundancy and makes the database more efficient and reliable.

# Explain how you can optimize the performance of SQL queries on this dataset.

There are several ways to optimize SQL queries for faster performance, a few are listed below:
1. Minimize the use of wildcard characters.
2. Increase Query Performance with Indexes.
3. Use appropriate data types.
4. Avoid subqueries.
5. Use LIMIT or TOP to limit the number of rows returned.
6. Avoid using SELECT * .
7. Use GROUP BY to group data.
8. Monitor query performance.